RECOMP II USERS' PROGRAM NO. 1093

PROGRAM TITLE:          STEPWISE MULTIPLE LINEAR REGRESSION
                        AND CORRELATION ANALYSIS

PROGRAM CLASSIFICATION: General Topics

AUTHOR:                 S. J. Singer
                        Aurora Gasoline Company

PURPOSE:                This program allows data for correlation
                        to be prepared on the versatape in a
                        simple format. It allows any variables
                        from the data tape to be selected for
                        correlation, and computes any desired
                        functions of the selected variables. The
                        program allows any variable to be selected
                        as a dependent variable and then computes
                        the statistically significant multiple
                        correlation between the dependent variable
                        and the rest of the variables selected.
                        As options the program will compute means,
                        standard deviations and partial correla-
                        tion coefficients.

DATE:                   1 August 1961

REVISED:                5 September 1961

Published by

RECOMP Users' Library

at

AUTONETICS INDUSTRIAL PRODUCTS

A DIVISION OF NORTH AMERICAN AVIATION, INC.

3400 E. 70 Street, Long Beach 5, Calif.

## DATA PREPARATION

## MULTIPLE REGRESSION PROGRAM

Introduction:

The data preparation program for the Multiple Regression Program serves three principle functions:

(1) It allows the original data to be prepared off line in a simple fixed point format.

(2) It permits data to be selected from the preliminary data tape in any desired fashion. Any variable on the preliminary data tape may be the dependent variable; many more variables may be put on the tape than are to be used in a single regression. A great deal of flexibility in data preparation is thus available.

(3) Perhaps most important, almost any function of the variable or combinations of variables on the preliminary data tape can be easily computed. Simple sums or products of variables, as well as logs, exponentials, sines and cosines or combinations of these functions can be computed.

The program uses a language similar to RAFT for formula compilation, however, the compiler produces a fairly efficient machine language program for each function typed in; thus, running time is kept to a minimum.

## MULTIPLE REGRESSION OPERATING INSTRUCTIONS

I.  **Input Routine General**

    A.   Punch   data on versa-tape using decimal interger cartridge.

    B.   All   data must be scaled $10^{-5}$  i.e. (XXXXX.XXXXX)

    C.   Each experiment should begin with L 3000

    D.   Each experiment should end with L 0000 "S"

    E.   A maximum of 32 data points including the dependent variable but not including the  constant may be included in each experiment.

    F.   There is no limit to the number of experiments.

    G.   The variables  are numbered and  referred to as 01 - 32 respectively.

II.  **Input of the Headings**

    A.   Three numbers must be typed in for each problem. These are:

        1.  The number of Independent variables (including the constant) N

        2.  The number of experiments M

        3.  An Identification number

    B.   These numbers  are entered from the typewriter as fixed point intergers.

    C.   The identification number may  be either positive or negative.  It controls the treatment of the constant in the regression equation  (See Regression Program Operating Instructions IV)

III.  **Equation Input**

    A.   Input is similar to RAFT except the beginning symbol for FCA (,) may be omitted.

    B.   Each equation should be ended with a TAB

    C.   Spaces and Figure shifts are ignored in equations

    D.   Three  temporary storage locations 90, 91 and 99 are available plus a pseudo accumulator 99.

    E.   The following functions are available as one digit codes:

        1.   FCA   ,    (or separator)

        2.   FAD   +

        3.   FSB   -

        4.   FMP   .

        5.   FDV   /

        6.   FSQ   ?

        7.   PWR   !

        8.   FST   :

F. Several subroutines are also available

    1. L/S I  invert (reciprocal)
    2. L/S E  $e^x$
    3. L/S L  Ln (base e)
    4. L/S G  Log (base 10)
    5. L/S X  $10^x$
    6. L/S S  Sine
    7. L/S C  Cosine

G. Constants may be entered by enclosing them in parentheses -- e.g. (+2.5) etc.

H. If a subroutine I, E, etc. is used it must either be separated from the rest of the equation by an arithmetic code, +, - etc. or by a comma.

I. The maximum number of characters in an equation is 256. (Figure Shifts & Spaces not included)

## IV. To Operate Program

A. Fill program tape - All tabs must be set

B. Press start "1" with sense switch "B" up
    (1) Computer types <u>Variables</u>:
        Type in No. of independent variables C/R
    (2) Computer types <u>Observations</u>:
        Type in No. of experiments or observations C/R
    (3) Computer types IDENT No.:
        Type an identification number C/R

C. In the above steps no comma sign or period need be typed.

D. The computer will type EQUATIONS then Y 00 followed by N X equations. For N equals 4,
The computer will type:  (Note X01 always equals one)
      Y 00
      X 01 (+1.0)
      X 02
      X 03
      X 04

E. As each equation number is typed, the appropriate equation should be entered, followed by a tab.
    (1) Examples X02 etc. =
      (a)  03 times 05                03.05 T/B
      (b)  Ln 3.5 div. by 01        L/S L(3.5)/01 T/B
      (c)  01 div. by (02 plus 03)  02+03, L/SI99.01 T/B
      (d)  Ln (01-05) div. by (01-05)  01-05:90,L/S L99/90 T/B
      (e)  $Log_{10} (1+e^{04})$         L/SE04+(1), L/S G99 T/B
      (f)  $Ln 12 \sqrt{e^{05}}$         L/S E05? 99:91, L/S L(12).9
                                        T/B
      (g)  05                       05 T/B

    (2) Typing line feed will erase the current equation or start "1" with SSB down will also erase an error.

F.  After the final equation has been entered the computer
    will type PUNCH DATA.  The preliminary data type punched
    on the versa tape should be placed in the photo reader.

G.  Press start 2 with sense switch G up.  Computer will
    read preliminary data tape, compute the variables, and
    punch a floating point tape in proper format for
    regression program.

H.  After punching is finished END will be typed out.
    Remove tape from punch and proceed to Regression Program.

I.  If it is desired to list the floating point data tape,
    load it in photo reader and press start 2 with sense
    switch "C" down.

J.  If it is desired to list the preliminary fixed point
    data tape, load it in the photo reader and press start.
    2 with sense switch "B" down.

    (1)  Computer types Input Lister, then Variables:
             Type in No. of independent variables C/R

    (2)  Computer types Observations:
             Type in No. of experiments or observations C/R

## Multiple Regression Program Description

### I. Introduction

Multiple Regression is a familiar statistical technique with a variety of uses. It may be used in the usual fashion to find the correlation between a single dependent variable and several independent variables or simply as a curve fitting program. Mathematically the problem may be stated as follows:

Given $t$ sets of observations having one dependent variable $Y$ and $n$ independent variables $X_1$, $X_2$,...$X_n$ fine the coefficients $c_0$, $c_1$, $c_2$, ...$c_n$ in the equation

$$Y = c_0 + c_1X_1 + c_2X_2 + ... + c_nX_n$$

which reduce the sum of squares of $Y$ the most. That is

$$\sum_{t=1}^{t} (Y \text{ observed} - Y \text{ predicted})^2 = \text{minimum}$$

### II. The method

The method used is essentially the stepwise multiple regression procedure given by Efroymson[1]. The computational procedure is as follows. The data is read in one experiment at a time and the following matrix is computed, each element being a sum of $t$ cross products.

Let the Matrix M

| i \ j | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| 0 | $Y^2$ | $X_1Y$ | $X_2Y$ | $X_3Y$ |
| 1 | $X_1Y$ | $X_1X_1$ | $X_2X_1$ | $X_3X_1$ |
| 2 | $X_2Y$ | $X_1X_2$ | $X_2X_2$ | $X_3X_2$ |
| 3 | $X_3Y$ | $X_1X_3$ | $X_2X_3$ | $X_3X_3$ |

be represented

$$\begin{array}{c|c} Y & c^T \\ \hline C & X \end{array}$$

Then inverting the matrix X by a series of linear transformations will transform the column C into the multiple regression coefficients.

This inversion process is carried out stepwise by successively pivoting on diagonal elements of the X matrix. Each step constitutes one iteration. Pivoting on the diagonal element corresponding to a variable not in the solution brings this variable into the solution. Pivoting on the diagonal element corresponding to a variable already in the solution causes this variable to be dropped.

At each iteration an F test is made to determine which variable to add to most greatly improve "the goodness of fit." After several variables

have been added one of the previously added variables may become statistically insignificant and will then be dropped. This stepwise procedure assures that only significant variables are included in the final solution.

An additional algorithm is included in the program as an option, which will cause all the variables to be added before any statistical testing takes place.

Unlike many regression programs this regression is not automatically fitted about the mean. Instead an extra column is carried in the matrix for the constant. This permits handling the constant just like any other independent variable, adding or dropping it at will.

III. Symbols and Definitions

$a_{ij}$ = any element in the matrix M

$a_{kk}$ = a pivot element in the matrix M

F to enter = F value for entering a variable

F to drop = F value for dropping a variable

i = row number of the matrix

j = column number of the matrix

k = index of a pivot element in the matrix

M = the matrix described in section II

N = number of independent + dependent variables

n = number of experiments

$s_y$ = standard error of dependent variable

V minus = variance increase caused by dropping the least significant variable in the solution

V plus = variance decrease caused by adding the most significant variable not in the solution

$x_i$ = the ith variable

$\bar{x}_i$ = the mean of the ith variable

Y = element $a_{oo}$ in matrix M

$\phi$ = degrees of freedom

1. Mean of variable $= \dfrac{\sum x_i}{n}$

2. Standard error of a variable $= \sqrt{\dfrac{\sum (x_i - \bar{x}_i)^2}{n-1}}$

3. Partial correlation coefficient for variables $x_i$ and $x_j$ $= \dfrac{\sum (x_i - \bar{x}_i)(x_j - \bar{x}_j)}{\sqrt{\sum (x_i - x_i)^2 \ \sum (x_j - \bar{x}_j)^2}}$

4. General algorithm for pivoting on element $a_{kk}$

$$\text{The new } a_{ij} = \begin{cases} \dfrac{a_{ij}a_{kk} - a_{ik}a_{kj}}{a_{kk}} & \text{if } i \neq k; \ j \neq k \\[3mm] \dfrac{a_{kj}}{a_{kk}} & \text{if } i = k; \ j \neq k \\[3mm] \dfrac{-a_{ik}}{a_{kk}} & \text{if } i \neq k; \ j = k \\[3mm] \dfrac{1}{a_{kk}} & \text{if } i = j = k \end{cases}$$

5. Variance change caused by pivoting on element $a_{kk}$ $= V = \dfrac{a_{ok} \ a_{ko}}{a_{kk}}$

6. Degrees of freedom $= \phi = N - $ number of independent variables in solution

7. F test for dropping a variable $= F \text{ to drop } \dfrac{|V \text{ minus}| \ \phi}{Y}$

8. F test for entering a variable $= F \text{ to enter } \dfrac{(V \text{ plus})(\phi)}{Y - V \text{ plus}}$

9. Standard error of dependent variable $= s_y = \sqrt{\dfrac{Y}{\phi}}$

10. Confidence limit of the regression coefficient of variable $x_i$ $= t \ s_y \sqrt{a_{ii}}$

## IV. Program Description

The program occupies twenty channels of memory from 00 through 23. Channels 24 - 33 are available for more program. Channel 34 contains a t table for 95% confidence limits, channel 35 contains the names of the variables, 36 is used for temporary storage and the matrix M is stored in channels 37 through 76.

The computation is done floating point. Running time depends on the size of the problem and output desired. For a typical problem with 16 variables the time is about 20 seconds per experiment to compute the matrix and about 25 seconds per iteration.

It is recommended that the data be scaled so that the numbers are between 0.1 and 10. With this scaling a 16 variable problem has been run over 30 iterations with no noticable round off error to five digit accuracy.

The program contains only one error halt. If a pivot element becomes zero or negative the program will type out <u>Pivot Element Too Small</u> and halt.

## V.   References

1.  Efroymson, M.A.   Stepwise Procedure for Calculation of Multiple
    Regression.   Presentation at Gordon Research Conference on Statistics.
    August 8-12, 1955.

2.  Duncan, Acheson J.   Quality Control and Industrial Statistics.
    Richard D. Irwin, Inc. 1953

3.  Ralston, Anthony & Wilf, Herbert S. Editors,   Mathematical Methods
    for Digital Computers, John Wiley & Sons Inc. 1960.

4.  Mickley, Sherwood & Reed, Applied Mathematics in Chemical Engineering
    McGraw Hill Book Company, Inc. 1957.

# Multiple Regression General Flow Diagram

A

Type Iteration Heading

Compute Addr of
pivot element which
reduces $\Sigma Y^2$ max $P_a$
and Addr of pivot
element which increases
$\Sigma Y^2$ max $P_d$ and

Variances for both

(Y) — Test F to drop > $F_d$ — (N)

Set pivot Addr to
$P_d$

(Y) — Test F to add > $F_a$ — (N)

Set pivot Addr to
$P_a$

Pivot on Pivot element
Type stand error + Sum Sq.

Output Headings
degrees of freedom
initial + final st error
of Y and F ratio

Output variables in
Solution 95% confidence
+ Sum sq red
Output variables not in
solution + sum sq red
Halt

(START)

Type Load Data Tape
Halt

Compute + Output
Observed vs predicted
Results + st error

# REGRESSION PROGRAM

## Multiple Regression Operating Instructions

I. Tape Format

    A. Tape format is Alphabetic, with the experiments in order.

    B. Each block begins with an "F" and ends with an "S".

    C. A leading block must head each problem tape. This block contains 3 words in the following order (fixed point @ B39):

        1. The No. of independent variables including the constant

        2. The No. of experiments

        3. An Identification No.

    D. Each experiment must have the variables on tape in the following order (floating point)

    E. A total of 31 independent variables including the constant (one) is the maximum allowable.

    F. There is no limit to the number of experiments.

II. The Sense Switch Settings

    A. Sense Switch B on causes the means standard deviations and partial correlation coefficients to be typed out.

    B. Sense Switch C on causes the entire matrix to be inverted first before any statistical testing takes place.

    C. Sense Switch D on causes all output to be in Floating Point; otherwise the output is Fixed Point 10 digits scaled $10^{-5}$.

III. The Start Switches

    A. Start one clears memory, sets the F levels to constant values (around 4.5) and starts reading the input tape.

    B. Start two allows new F values to be entered from the typewriter and begin iterating.

    C. Start three same as start one except computer halts to allow F values to be entered from typewriter.

IV. Control of the Constant

A. The program is set up with three possibilities concerning the constant (variable X01). These are controlled by the sign of the ID Number.

1. If the ID No. is positive the computer will add or drop the constant according to statistical tests.

2. If the ID No. is negative the choice of whether or not the constant is to be added is up to the operator. If the ID No. is negative the computer will halt after reading in the data tape and type:

SET CONSTANT     Type , + 1. tab s or
                      , + 0. tab s

Typing a one will force the constant into the answer, typing a zero will leave it out.

3. The sign of the ID number may, of course, be changed at any time to obtain the desired result. The ID number is stored in location 7702 octal.

V. Forcing a Solution

A. It is sometimes desired to force a given solution to a problem. Provision to do this has been provided in the program.

B. To force a given solution

1. Transfer to Location 0000.0 (Type L 00000 Enter on console)

2. Computer will type ENTER REQUIRED SOLUTION, then CHANGE VARIABLE NO.  Type , + XX. Tab s

(a) If variable XX is in the solution, it will be forced out

(b) If it is not in the solution, it will be forced in

3. Computer will perform the desired operation and return to the point where CHANGE VARIABLE NO. is typed out.

4. Typing a zero for the variable (, + 0. Tab s) will cause the computer to type out MULTIPLE REGRESSION REQUIRED SOLUTION then enter the normal output routine.

VI. The Output

A. The program always outputs the following information.

1. The Problem Identification No.

VI. The Output (continued)

    2. As each variable is added or dropped the computer outputs:

       (a) The Name of the variable

       (b) The Sum of Squares of Y before adding this variable

       (c) The Standard Error of Y

    3. When the Optimum solution is reached, the Computer types out:

       (a) The names of the Variables in the solution

       (b) The initial and final Sums of Squares and Standard Errors of the dependent variable Y

       (c) The F ratio for the regression

    4. The Type-out continues with:

       (a) The names of the Variables in the solution

       (b) The regression coefficients of those variables

       (c) The 95% confidence limits of the coefficients

       (d) The amount the sum of squares of Y would be increased if the variable were added to the solution.

B. As optional additional information the following can be obtained (see sense switch settings and operating instructions)

    1. The means and standard deviations from the mean of each variable.

    2. The partial correlation coefficients of the variables.

    3. The observed and predicted value of the dependent variable for each experiment and the deviation, and a recomputed value for the standard error of Y calculated from the original data.

## VII. To Operate Program

A.  Fill Program - All tabs must be set

B.  Load Floating Point data tape in photoreader and press Start 1 or 3

    1.  If Start 1 is pressed, the computer will read the leading block type out the Id No., read the rest of the data and begin iterating

    2.  If Start 3 is pressed, the computer will halt and type

|  |  |  |  |
|---|---|---|---|
| F to ENTER | Type , + X.X | Tab | S |
| F to DROP | Type , + X.X | Tab | S |

    Computation will continue as in the case where Start 1 was pressed.

    3.  When entering F levels the comma (,) and sign (+) must be typed. The F values are entered as mixed numbers; up to 10 digits may be entered both before and after the decimal point.

C.  Computer will add and drop variables as determined by the sense switch settings and F levels typing out intermediate results as it does so.

D.  After an optimum solution has been reached, the computer will output the answers and halt

E.  To run observed vs Predicted Results press Start

    1.  Computer will type Load Data Tape and halt.

    2.  Load the data tape in the photoreader and press Start

    3.  Computer will read tape & output the observed predicted values and deviation of the dependent variable.

    4.  When all variables have been computed the computer will type out the recomputed standard error.

## VIII. Example Problem

The following is a sample problem taken from Duncan[2], Quality Control and Industrial Statistics, which illustrates several of the options in the program.

The problem has one dependent variable, four independent variables, including the constant and twenty experiments. A regression was run using F levels of 4.0. This resulted in three of the variables entering the solution. This solution is shown on Page 5 of the sample problem type-out. The remaining variable was then forced into the solution. This solution is shown on Page 6 of the sample problem type-out.

# MULTIPLE REGRESSION SAMPLE PROBLEM

| Experiment No. | Y00 | X02 | X03 | X04 |
|---|---|---|---|---|
| 1. | 9.9 | 8.5 | 7.6 | 4.4 |
| 2. | 9.3 | 8.2 | 7.8 | 4.2 |
| 3. | 9.9 | 7.5 | 7.3 | 4.2 |
| 4. | 9.7 | 7.4 | 7.2 | 4.4 |
| 5. | 9.0 | 7.6 | 7.3 | 4.3 |
| 6. | 9.6 | 7.4 | 6.9 | 4.6 |
| 7. | 9.3 | 7.3 | 6.9 | 4.6 |
| 8. | 13.0 | 9.6 | 8.0 | 3.6 |
| 9. | 11.8 | 9.3 | 7.8 | 3.6 |
| 10. | 8.8 | 7.0 | 7.3 | 3.7 |
| 11. | 8.9 | 8.2 | 7.1 | 4.6 |
| 12. | 9.3 | 8.0 | 7.2 | 4.5 |
| 13. | 9.4 | 7.7 | 7.6 | 4.2 |
| 14. | 7.5 | 6.7 | 7.6 | 5.0 |
| 15. | 8.4 | 8.2 | 7.0 | 4.8 |
| 16. | 9.1 | 7.6 | 7.6 | 4.1 |
| 17. | 10.0 | 7.4 | 7.8 | 3.1 |
| 18. | 9.8 | 7.1 | 8.0 | 2.9 |
| 19. | 10.1 | 7.0 | 8.3 | 3.9 |
| 20. | 8.0 | 6.4 | 7.9 | 3.8 |

Y00 = dependent variable

X02 - X04 = independent variables

# MULTIPLE REGRESSION SAMPLE PROBLEMS
## DATA PREPARATION

**Data Preparation      Start 1**

VARIABLES:      4
OBSERVATIONS: 20

IDENT NO:      2

EQUATIONS:

| | |
|---|---|
| Y00 | 01 |
| X01 | (+1.0) |
| X02 | 02 |
| X03 | 03 |
| X04 | 04 |

PUNCH DATA                    **Load Preliminary   Data Tape**

**Listing of Compiled Floating Point**

| | | |
|---|---|---|
| N 001 | | |
| 00 | 9.90000 | **Tape      Start 2 with SSC down.** |
| 01 | 1.00000 | |
| 02 | 8.50000 | |
| 03 | 7.60000 | |
| 04 | 4.40000 | |
| | | |
| N 002 | | |
| 00 | 9.30000 | |
| 01 | 1.00000 | |
| 02 | 8.20000 | |
| 03 | 7.80000 | |
| 04 | 4.20000 | |
| | | |
| N 003 | | |
| 00 | 9.90000 | |
| 01 | 1.00000 | |
| 02 | 7.50000 | |
| 03 | 7.30000 | |
| 04 | 4.20000 | |
| | | |
| N 004 | | |
| 00 | 9.70000 | |
| 01 | 1.00000 | |
| 02 | 7.40000 | |
| 03 | 7.20000 | |
| 04 | 4.40000 | |

```
N 005
  00        9.00000
  01        1.00000
  02        7.60000
  03        7.30000
  04        4.30000

N 006
  00        9.60000
  01        1.00000
  02        7.40000
  03        6.90000
  04        4.60000

N 007
  00        9.30000
  01        1.00000
  02        7.30000
  03        6.90000
  04        4.60000

N 008
  00       13.00000
  01        1.00000
  02        9.60000
  03        8.00000
  04        3.60000

N 009
  00       11.80000
  01        1.00000
  02        9.30000
  03        7.80000
  04        3.60000

N 010
  00        8.80000
  01        1.00000
  02        7.00000
  03        7.30000
  04        3.70000

N 011
  00        8.90000
  01        1.00000
  02        8.20000
  03        7.10000
  04        4.60000

N 012
  00        9.30000
  01        1.00000
  02        8.00000
  03        7.20000
  04        4.50000
```

N 013
| | |
|---|---|
| 00 | 9.40000 |
| 01 | 1.00000 |
| 02 | 7.70000 |
| 03 | 7.60000 |
| 04 | 4.20000 |

N 014
| | |
|---|---|
| 00 | 7.50000 |
| 01 | 1.00000 |
| 02 | 6.70000 |
| 03 | 7.60000 |
| 04 | 5.00000 |

N 015
| | |
|---|---|
| 00 | 8.40000 |
| 01 | 1.00000 |
| 02 | 8.20000 |
| 03 | 7.00000 |
| 04 | 4.80000 |

N 016
| | |
|---|---|
| 00 | 9.10000 |
| 01 | 1.00000 |
| 02 | 7.60000 |
| 03 | 7.60000 |
| 04 | 4.10000 |

N 017
| | |
|---|---|
| 00 | 10.00000 |
| 01 | 1.00000 |
| 02 | 7.40000 |
| 03 | 7.80000 |
| 04 | 3.10000 |

N 018
| | |
|---|---|
| 00 | 9.80000 |
| 01 | 1.00000 |
| 02 | 7.10000 |
| 03 | 8.00000 |
| 04 | 2.90000 |

N 019
| | |
|---|---|
| 00 | 10.10000 |
| 01 | 1.00000 |
| 02 | 7.00000 |
| 03 | 8.30000 |
| 04 | 3.90000 |

N 020
| | |
|---|---|
| 00 | 8.00000 |
| 01 | 1.00000 |
| 02 | 6.40000 |
| 03 | 7.90000 |
| 04 | 3.80000 |

END

| F TO ENTER | ,+4.0  s | Fill Regression Program |
|---|---|---|
| F TO DROP | ,+4.0  s | Start 3 with Sense Switch B Down |

| MULTIPLE REGRESSION PROBLEM NUMBER | 2.00000 |
|---|---|

## MEANS AND STANDARD DEVIATIONS

| | MEANS | STAND. DEV. |
|---|---|---|
| YOO | 9.54000 | 1.20149 |
| XO1 | 1.00000 | .00000 |
| XO2 | 7.70500 | .80032 |
| XO3 | 7.51000 | .39987 |
| XO4 | 4.12500 | .54664 |

## PARTIAL CORRELATION COEFFICIENTS

YOO VERSUS
XO1
| XO2 | .72940 |
|---|---|
| XO3 | .36721 |
| XO4 - | .48723 |

.XO1 VERSUS
XO2
XO3
XO4

XO2 VERSUS
| XO3 | .04096 |
|---|---|
| XO4 - | .03519 |

XO3 VERSUS
| XO4 - | .66577 |
|---|---|

| | VAR | SUM SQ | ST ERROR |
|---|---|---|---|
| ADD | XO2 | 1847.66000 | 9.61161 |
| ADD | XO4 | 13.08236 | .82979 |
| ADD | XO1 | 10.64652 | .76908 |

## MULTIPLE REGRESSION OPTIMUM SOLUTION

DEGREES OF FREEDOM     17.00000

| | ST ERROR | SUM SQ |
|---|---|---|
| INITIAL | 9.61161 | 1847.66000 |
| FINAL | .64102 | 6.98526 |
| F RATIO | 3.51324 | |

## VARIABLES IN SOLUTION

| NAME | COEFFICIENT | CONFIDENCE | SUM SQ RED |
|------|-------------|------------|------------|
| X01  | 5.48080     | 3.87416    | 3.66127    |
| X02  | 1.07062     | .38796     | 13.93171   |
| X04 −| 1.01574     | .56799     | 5.85034    |

## VARIABLES NOT IN SOLUTION

| NAME | SUM SQ RED |
|------|------------|
| X03  | .04536     |

LOAD DATA TAPE

### OBSERVED VS PREDICTED RESULTS

| EXP  | OBSERVED | PREDICTED | DEVIATION |
|------|----------|-----------|-----------|
| 0001 | 9.90000  | 10.11182  | − .21182  |
| 0002 | 9.30000  | 9.99378   | − .69378  |
| 0003 | 9.90000  | 9.24435   | .65566    |
| 0004 | 9.70000  | 8.93414   | .76587    |
| 0005 | 9.00000  | 9.24983   | − .24983  |
| 0006 | 9.60000  | 8.73099   | .86902    |
| 0007 | 9.30000  | 8.62393   | .67608    |
| 0008 | 13.00000 | 12.10209  | .89792    |
| 0009 | 11.80000 | 11.78090  | .01910    |
| 0010 | 8.80000  | 9.21691   | − .41691  |
| 0011 | 8.90000  | 9.58749   | − .68749  |
| 0012 | 9.30000  | 9.47494   | − .17494  |
| 0013 | 9.40000  | 9.45847   | − .05847  |
| 0014 | 7.50000  | 7.57526   | − .07526  |
| 0015 | 8.40000  | 9.38434   | − .98434  |
| 0016 | 9.10000  | 9.45298   | − .35298  |
| 0017 | 10.00000 | 10.25459  | − .25459  |
| 0018 | 9.80000  | 10.13656  | − .33656  |
| 0019 | 10.10000 | 9.01376   | 1.08625   |
| 0020 | 8.00000  | 8.47296   | − .47296  |

ST ERROR      .64102

ENTER REQUIRED SOLUTION

CHANGE VARIABLE NO ,+3. s

| | VAR | SUM SQ | ST ERROR |
|---|---|---|---|
| ADD | X03 | 6.98526 | .64102 |

CHANGE VARIABLE NO ,+0. s

MULTIPLE REGRESSION REQUIRED SOLUTION

DEGREES OF FREEDOM 16.00000

| | ST ERROR | SUM SQ |
|---|---|---|
| INITIAL | 9.61161 | 1847.66000 |
| FINAL | .65860 | 6.93990 |
| F RATIO | 3.32807 | |

VARIABLES IN SOLUTION

| NAME | COEFFICIENT | CONFIDENCE | SUM SQ RED |
|---|---|---|---|
| X01 | 3.93295 | 10.90709 | .25347 |
| X02 | 1.06919 | .40060 | 13.88670 |
| X03 | .16381 | 1.07386 | .04536 |
| X04 - | .93604 | .78536 | 2.76922 |

VARIABLES NOT IN SOLUTION

| NAME | SUM SQ RED |
|---|---|

LOAD DATA TAPE